

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Unpacking the functional heterogeneity of the Emotional Face Matching Task: a normative modelling approach

Hannah S. Savage^{1,2}, Peter C. R. Mulders^{1,3}, Philip F. P. van Eijndhoven^{1,3}, Jasper van Oort^{1,3}, Indira Tendolkar^{1,3}, Janna N. Vrijsen^{1,3,4}, Christian F. Beckmann^{1,2,5}, Andre F. Marquand^{1,2}

¹Donders Institute of Brain, Cognition and Behaviour, Radboud University, Nijmegen, The Netherlands

²Department of Cognitive Neuroscience, Radboud University Medical Centre, Nijmegen, The Netherlands

³Department of Psychiatry, Radboud University Medical Centre, Nijmegen, The Netherlands

⁴Depression Expertise Centre, Pro Persona Mental Health Care, Nijmegen, The Netherlands

⁵Centre for Functional MRI of the Brain (FMRIB), Nuffield Department of Clinical Neurosciences, Wellcome Centre for Integrative Neuroimaging, University of Oxford, Oxford, UK

Word count main text: 4745

Figures: 6

Tables: 1

Corresponding author:

Hannah S. Savage | Andre F. Marquand

Kappitelsweg 29,

Nijmegen, 6525EN

The Netherlands

E: hannah.savage@donders.ru.nl | andre.marquand@donders.ru.nl

Author contribution statement:

H.S.S – conceptualization, formal analysis, data curation, writing – original draft, writing – review & editing, visualisation. P.C.R.M – formal analysis, data curation, resources, writing – review & editing. P.F.P.E – resources, writing – review & editing, funding acquisition. J.O. – investigation, data curation, writing – review & editing. I.T. – resources, writing – review & editing. J.N.V. – resources, writing – review & editing. C.F.B. – resources, funding acquisition, writing – review & editing. A.F.M – conceptualization, supervision, funding acquisition, formal analysis, writing – original draft, writing – review & editing.

32 **Abstract:**

33 Functional neuroimaging has contributed substantially to understanding brain function but is dominated by
34 group analyses that index only a fraction of the variation in these data. It is increasingly clear that parsing
35 the underlying heterogeneity is crucial to understand individual differences and the impact of different task
36 manipulations. We estimate large-scale (N=7641) normative models of task-evoked activation during the
37 Emotional Face Matching Task, which enables us to bind heterogeneous datasets to a common reference
38 and dissect heterogeneity underlying group-level analyses. We apply this model to a heterogeneous patient
39 cohort, to map individual differences between patients with one or more mental health diagnoses relative
40 to the reference cohort and determine multivariate associations with transdiagnostic symptom domains. For
41 the face>shapes contrast, patients have a higher frequency of extreme deviations with unique spatial
42 distributions depending on diagnosis. In contrast, normative models for faces>baseline have greater
43 predictive value for individuals' transdiagnostic functioning.

44 **Introduction:**

45 Task-based functional neuroimaging (functional magnetic resonance imaging; fMRI) has been widely
46 applied in foundational and clinical neuropsychology to characterise neural processes that underpin a
47 behaviour or process of interest. The typical approach in such studies is based on comparing mean
48 differences in the magnitude and location of activation (measured by changes in BOLD signal), which has
49 helped us to understand how these processes may differ between groups defined by biological and
50 sociocultural factors, psychopathologies, or therapeutic interventions. The majority of prior research has
51 reported group-level summary statistics, which inform us of those regions most consistently activated
52 across participants/groups during task conditions. This method assumes that the neural mechanisms
53 facilitating the process of interest are consistent across individuals within and between groups. This
54 assumption enables our understanding to reach only so far as ‘the average brain’ of an ‘average control’,
55 or ‘average patient’.

56 In order to better understand how the brain relates to behaviour it is essential to move our focus
57 from the group-level to studying individual differences and consider the neural activation of these processes
58 within the context of multiple sources of heterogeneity. For example: (i) natural variation within the general
59 population, including potentially heterogenous yet functionally convergent processes, and (ii) heterogeneity
60 within groups of interest, such as within mental health diagnoses. Furthermore, when comparing between
61 independent studies, the influence of task design (i.e. small modifications to an original task) and acquisition
62 parameters should also be considered but are seldom investigated.

63 One approach that can provide insight into individual differences is normative modelling^{1,2}. The
64 normative modelling framework provides statistical inference at the level of each subject with respect to an
65 expected pattern across the population, highlighting variation within populations in terms of the mapping
66 between biological variables and other measures of interest. This framework has previously been employed
67 by our group and others to dissect structural variation within large healthy populations³ and clinical
68 psychiatric populations (e.g. in autism⁴⁻⁶, schizophrenia and bipolar disorder⁷), and in relation to dimensions
69 of psychopathology⁸. Applying this method to task-based fMRI data we will be able to characterize how
70 functional activity within each voxel or ROI in the brain differs between individuals, and hence show with
71 greater nuance the range of task-evoked activation within the general population². Further, applying this
72 model to patients with a current diagnosis (mood and anxiety disorders, autism spectrum disorders (ASD)
73 and/or attention deficit hyperactivity disorder (ADHD)) we will be able to map differences in these individual
74 participants with respect to the reference cohort. This may reveal unique clusters of deviation patterns,
75 within and/or across diagnostic categories.

76 In this study, we use the Emotional Face Matching Task (EFMT) to demonstrate the potential of
77 the normative modelling method to identify individual differences in task-based fMRI. The EFMT, also
78 commonly referred to as the ‘Hariri task’, has been used in over 250 fMRI studies since it was most notably
79 introduced in 2002^{9,10}. This task asks participants to match one of two images that are simultaneously

80 presented at the bottom of the screen, to a third target image displayed at the top of the screen; participants
81 match images of facial configurations consistent with the common view of prototypic facial expressions,
82 most frequently of fear or anger, or similarly positioned geometric shapes. Matching faces, as compared to
83 matching shapes, evokes explicit and/or implicit emotional face processing, which has previously been
84 shown to engage the amygdala, fusiform face area, anterior insula cortex, the pregenual and dorsal anterior
85 cingulate cortex, the dorsomedial and dorsolateral prefrontal cortex, and visual input areas. Previous work
86 has related activity to biological and demographic variables, and compared between many different clinical
87 groups and developmental spectrums.

88 Due to its experimental simplicity and focus on subcortical circuitry relevant to brain disorders, the
89 EFMT has been implemented in a number of large-scale neuroimaging initiatives including the UK
90 Biobank¹¹, the Human Connectome Project (HCP)^{12,13}, HCP Development¹⁴, the Amsterdam Open MRI
91 Collection Population Imaging of Psychology (AOMIC PIOP2)¹⁵, and the Duke Neurogenetics Study (DNS).
92 We take advantage of these large open-access/shared datasets to collate a large representative sample of
93 over 7500 participants from six sites to first (1) build reference normative models that highlight the natural
94 variation of functional activity evoked by the EFTM [as measured by the task contrasts faces > shapes and
95 faces>baseline], and (2) determine how the model's prediction relates to age, sex, and variations in task
96 design. We then apply these models to over 200 participants with a current mental health condition or who
97 are neurodivergent from the MIND-Set cohort (Measuring Integrated Novel Dimensions in
98 neurodevelopmental and stress-related psychiatric disorder)¹⁶, to (3) map deviations in patients with a
99 current diagnosis (mood and anxiety disorders, ASD and/or ADHD) relative from the reference cohort, both
100 at the group level and at the level of the individual. We show that despite the ostensible simplicity of this
101 task and robust group effects, there is considerable inter-individual heterogeneity in the nature of the elicited
102 activation patterns and that such differences are both highly interpretable and predict cross-domain
103 symptomatology in a naturalistic clinical cohort.

104

105 **Results:**

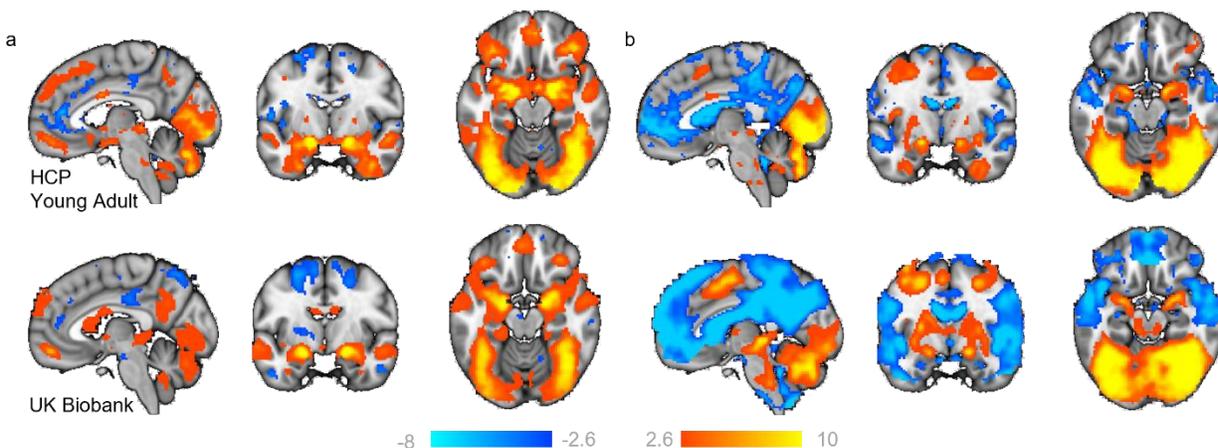
106 ***Group level comparisons show consistent effects across cohorts***

107 First, we performed a classical group comparison to provide a reference against which to understand the
108 inter-individual differences in subsequent analyses. To achieve this, we randomly selected 100 random
109 individuals' FSL pre-processed data into fixed-effects general linear models to create group level maps for
110 the faces>shapes (Fig. 1a) and faces>baseline (Fig. 1b) contrasts (see methods). This also served as a
111 sanity check to ensure the data was comparable to past literature. Overall, positive task effects (activations)
112 for faces>shapes were found in the bilateral inferior and middle occipital lobe and the calcarine cortex (V1)
113 extending anterior-ventrally to the bilateral lingual and fusiform gyrus, and anterior-dorsally to the middle
114 and inferior temporal gyrus; the bilateral amygdala extending into the hippocampus; the bilateral temporal
115 pole; a dorsal region of the vmPFC; and the bilateral middle and inferior frontal gyrus. Task-related

116 deactivations were found across regions comprising the default mode network, including the anterior and
117 posterior cingulate cortex and precuneus, the precentral gyrus and supplementary motor area and the
118 inferior temporal lobe.

119

120



121

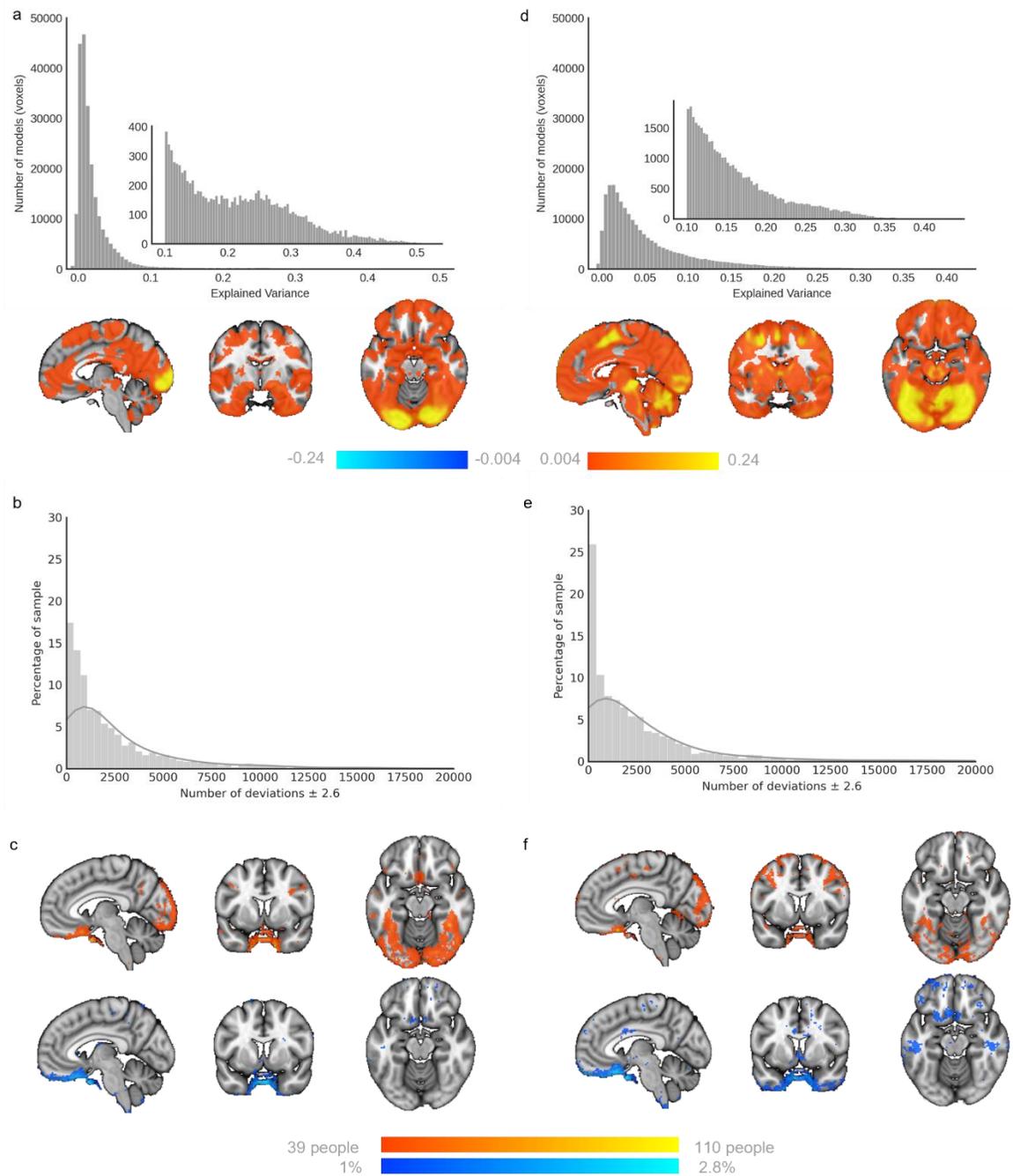
122 **Figure 1: Task evoked activation.** Two representative groups maps (from HCP Young Adult and UK Biobank), illustrating regions
123 where participants show greater BOLD signal (z-statistic maps, thresholded at ± 2.6) to (a) faces, as compared to shapes
124 (faces>shapes), and (b) faces, as compared to baseline (faces>baseline). $x,y,z = -4,-6,-16$.

125

126

127 **Fitting reference normative models for emotional face processing**

128 Next, we estimated normative models of EFMT-related BOLD activation for the face>shapes and
129 faces>baseline contrast using data from 7641 individuals across the lifespan. To achieve this, we split the
130 data into training ($n = 3877$) and test splits, stratified by site ($n = 3764$), then fit a Bayesian Linear
131 Regression model that predicted the single subject level activation for each voxel of the brain, as a function
132 of sex, age, and acquisition and task parameters (see methods). Explained variance in the test set was
133 good, especially in regions that showed activation at the group level (Fig. 1) including the occipital
134 lobe/visual cortex and the bilateral amygdala (faces>shapes: Fig. 2a; faces>baseline: Fig. 2d). As shown
135 in Supplementary Fig. 2a and 2c in most voxels the skew and kurtosis was acceptable (i.e. $-1 < \text{skew} < 1$
136 and kurtosis around zero). For a very small proportion of voxels this was not the case; the most ventral
137 region of the vmPFC (i.e. the bottom border of the brain) was the most negatively skewed which we interpret
138 to reflect the varying degrees of signal dropout, more so than biological variation. The few voxels with
139 positive kurtosis were spatially overlapping with regions that were negatively skewed, which likely reflects
140 the extended negative tails of the distributions in these voxels.



141

142 **Figure 2: Evaluation and deviation scores from the faces>shapes (left) and faces>baseline (right) normative models.**

143 Explained variance is high in the normative models, irrespective of whether they are built using the face>shapes contrast (a), or the

144 faces>baseline contrast ($x,y,z = -4,-6,-16$). (d). Histograms show the relative frequency of the total number of deviations that a

145 participant has for each model (b,e). Normative Probability Maps illustrate the percentage of participants of the total sample who had

146 positive (hot colours) or negative deviations (cool colours) $> \pm 2.6$ within each voxel, for the faces > shapes (c) and faces>baseline (f)

147 models. $x,y,z = -5,6,-15$.

148

149 ***Voxel-wise deviations show considerable inter-individual variability***

150 We then used these normative models to quantify the degree of inter-individual variability. To achieve this,
151 for each participant we created a thresholded normative probability map (NPM; deviation scores $>\pm 2.6$)
152 which indicates the difference between the predicted activation and true activation scaled by the prediction
153 variance, and therefore shows the voxels where that participant had greater or less activation than would
154 be expected by the normative models. Figures 2b and 2e show the frequency of the total number of
155 deviations that individuals had from the faces>shapes, and faces>baseline models, respectively. Within
156 each voxel, we then counted how many participants had positive or negative deviations ($>\pm 2.6$). The
157 resulting brain maps illustrate the variability in the magnitude of functional activation per voxel, across the
158 population for the two task contrasts (Fig. 2c + f). This shows that: (i) there is considerable inter-individual
159 variability underlying the mean effects and (ii) that the spatial distribution of individual deviations mostly
160 occurs within the task network. Every voxel of the brain had at least one subject with a deviation $>\pm 2.6$ (not
161 shown), although, as illustrated, there were regions including the medial occipital lobe extending to the
162 bilateral fusiform gyrus and inferior temporal lobe, the bilateral inferior frontal gyrus extending to the
163 precentral gyrus, and the posterior region of the vmPFC, wherein deviations were more frequently
164 observed. As there were minimal differences in the evaluation metrics between models built using either
165 contrast, and as the contrast faces>shapes is most commonly reported in prior literature, we use this as
166 our primary contrast for our further analysis of the reference model.

167

168 ***Separable effects of input variables on model predictions***

169 Next, we examine structure coefficients from our models to disentangle the effects of different input
170 variables. As shown in Figure 3, the direction of the relationship between input variables and the predicted
171 BOLD activation, and the fraction of the explained variability can be meaningfully separated for
172 interpretation. Some input variables, namely acquisition parameters, showed overlapping effects (with
173 sensical direction flips) likely due to their relatively high correlation and limited variability across sites. For
174 example, the number of target blocks, volumes acquired (not shown), use of multiband sequence (not
175 shown), and the length of the TR (not shown) all showed a similar relation to predicted activity.

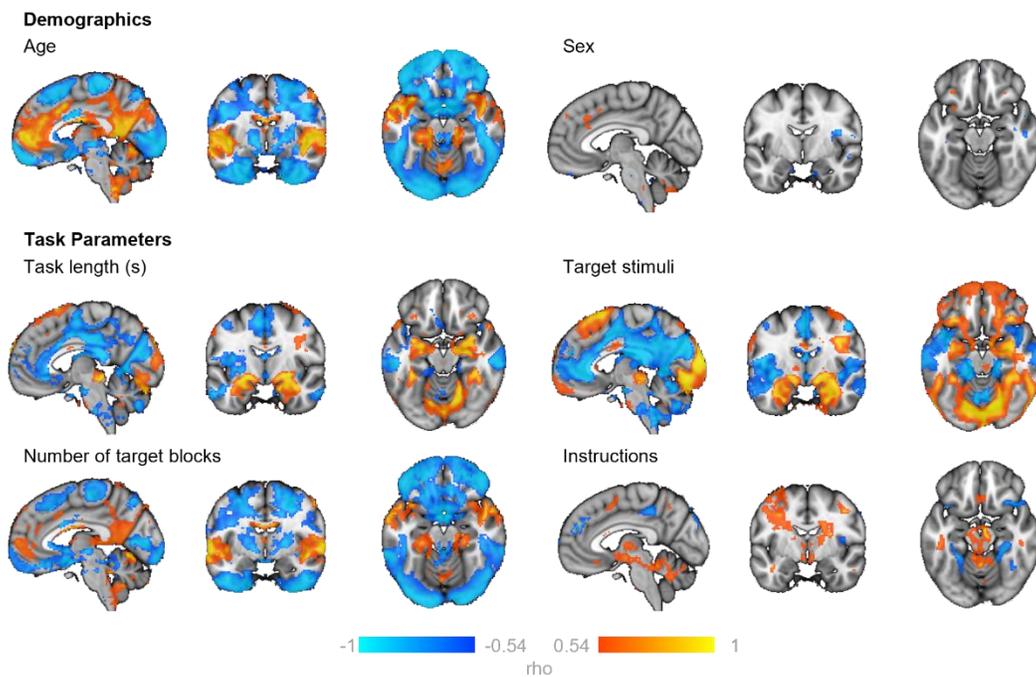
176 Increased age was related to decreased predicted activity across the peripheral/surface of the
177 brain, as well as regions surrounding the ventricles, and increased activity in midline regions of the default
178 mode network, the bilateral insula, the fusiform face area extending to the para-hippocampal gyrus and the
179 superior temporal gyrus. Predictions were only minimally influenced by sex, and the spatial mapping of this
180 relationship was broadly overlapping with that of intra-cranial volume (not shown).

181 We further illustrate the ability of this method to disentangle the influence of task design choices,
182 on predicted activation. For example, the influence of the matching rule and the stimuli presented. Being
183 told to match the emotional expression, as compared to matching the faces, related to increased predicted
184 BOLD activity within subcortical areas including the bilateral putamen, caudate body and medio-dorsal

185 thalamus. Attending to the emotional expression also predicted increased activity within the superior frontal
186 gyrus extending to the supplementary motor area, the posterior medial temporal gyrus the inferior temporal
187 gyrus, and the medial temporal pole. Conversely, when participants were asked to match faces, the model
188 predicted greater activation within the bilateral fusiform gyrus, the middle temporal gyrus, the superior
189 temporal pole, the dorsolateral prefrontal cortex, and a large area of the inferior parietal gyrus extending to
190 the supramarginal and angular gyrus. Additionally, when stimuli from the Ekman series were used the
191 model predicted greater activation within the bilateral inferior occipital gyrus and the calcarine cortex (V1),
192 the bilateral lingual and fusiform gyrus extending to the inferior temporal gyrus, as well as in the medial
193 cingulate cortex, an anterior region of the vmPFC, the superior medial prefrontal cortex, and subcortical
194 regions including the ventral posterior thalamus, the posterior putamen, para-hippocampus, hippocampus
195 and amygdala. Conversely, the use of the Nim-Stim Set stimuli related to greater activity within default
196 mode regions, including a large area of the ventromedial/medial prefrontal cortex, precuneus, cuneus, as
197 well as the supramarginal gyrus which extended medially to the anterior and posterior insula, which in turn
198 extended laterally to the superior and medial temporal gyri.

199

200



201

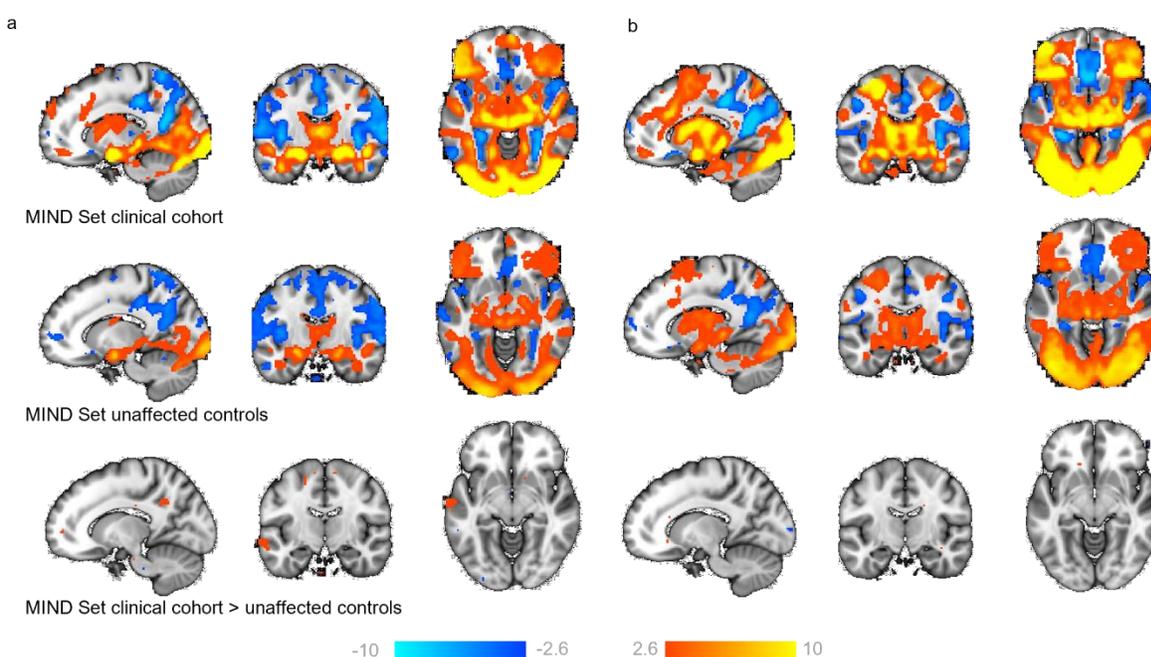
202 **Figure 3: The relationship between input variables and the predicted BOLD activation for faces > shapes.** Maps show the
203 correlation coefficients (ρ) thresholded by their respective coefficients of determination ($\rho^2 > 0.3$) of selected model input variables.
204 This can be interpreted as showing how predicted BOLD activation for the faces > shapes contrast relates to the input variables of the
205 normative models. Positive correlations (warm colours) indicate greater activation for higher values of the input variable and negative
206 correlations (cool colours) greater activation for lower values of the input variable (note that some variables are dummy coded, e.g.
207 target stimuli, instructions) $x, y, z = -4, -6, -16$.

208
209

210 **A traditional case-control comparison identifies few differences between patients and**
211 **controls**

212 We then performed a voxel-wise case-control comparison on the raw data in order to test for group level
213 differences between a heterogeneous patient cohort and matched controls from the naturalistic MIND-Set
214 sample. As evidenced in Table 1 (see Diagnoses), the naturalistic MIND-Set sample has many patients
215 with co-occurring and heterogenous mental health diagnosis, with or without neurodivergence, and is
216 therefore representative of diverse clinical populations. This analysis revealed very few differences between
217 the patient cohort, and unaffected controls for faces>shapes and faces>baseline (Fig. 4a and b – bottom
218 rows). More specifically, comparing patients' task activation (Fig. 4a – top row) to controls (Fig. 4a – middle
219 row) for the faces>shapes contrast showed patients had greater activation in the left temporal medial gyrus
220 and bilateral posterior cingulate cortex, as well as in small regions of the supplementary motor area, and
221 the genu of the anterior cingulate cortex (Fig. 4a – bottom row). There were negligible differences between
222 patients and controls for the faces>baseline contrast (Fig. 4b – bottom row).

223



224
225
226
227
228
229

Figure 4: General linear model results comparing patients to controls for the faces>shapes and faces>baseline contrasts. Maps show regions activated (warm colours) and deactivated (cool colours) for faces>shapes (a) and faces>baseline (b), for patients (top row) and unaffected controls (middle row) from the MIND Set cohort. (c) Regions where patients have more activation than controls (bottom row) (z-statistic maps, thresholded at $> \pm 2.6$). x,y,z = -14,-13,-9.

230

231 **Application of normative model to a naturalistic clinical sample**

232 Next, we aimed to relate the deviations from these normative models to psychopathology. To achieve this,
233 we evaluated the patient cohort with respect to the normative models estimated from the large reference
234 cohort. For the faces>shapes and faces>baseline models, the explained variance of the clinical test data
235 was quite low. This was expected given that this cohort is quite homogenous with respect to the covariates

236 included in the model (i.e., all subjects were scanned on the same scanner, using the same experimental
237 paradigm and had an age range considerably narrower than the reference cohort). This suggests that the
238 variance in BOLD signal was driven more by individual differences, as opposed to the variables included in
239 the model. The skew and kurtosis of the models were centred around zero. See Supplementary Figure 4a-
240 f for histograms of these evaluation metrics, and their respective illustration on the brain.

241

242 ***Frequency of deviations differentiates patients from unaffected controls***

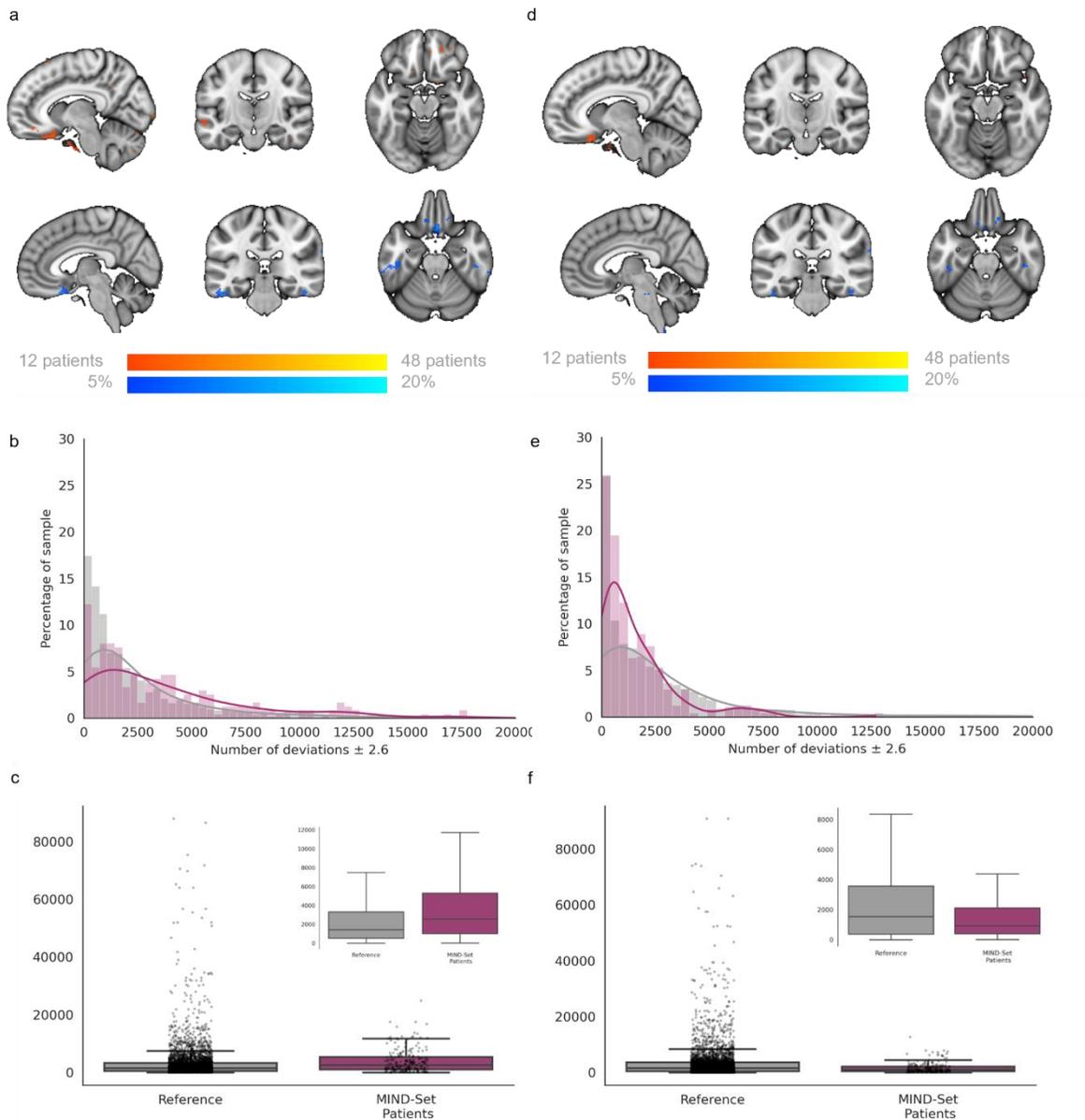
243 Next, we compared the frequency of extreme deviations (NPMs thresholded at $> \pm 2.6$), at the level of each
244 individual, between patients from the MIND-Set cohort and unaffected controls. For each model type
245 (faces>shapes: Fig. 5b,c; faces>baseline: Fig. 5e,f). MIND-Set patients had a greater frequency of
246 deviations relative to the reference cohort for the faces>shapes contrast (Mann-Whitney U test = 358986.5,
247 $p = 1.55^{-8}$; Fig. 5b). These deviations were most frequently identified in the lateral ventral prefrontal cortex,
248 and the bilateral medial and inferior temporal lobe (Fig. 5a). In contrast, for the faces>baseline contrast
249 individuals from the reference cohort had a greater frequency of deviations relative to MIND-Set patients
250 (Mann-Whitney U test = 509017.0, $p = 0.0007$; Fig. 5e). For this contrast, these deviations are, however,
251 strongly localised within the most ventral region of the vmPFC (Fig. 5d) which is well-known to be
252 problematic area for signal distortion artefacts in fMRI, therefore we do not interpret this difference as being
253 biologically meaningful.

254

255 ***Associations of patterns of deviation with cross-diagnostic symptom domains***

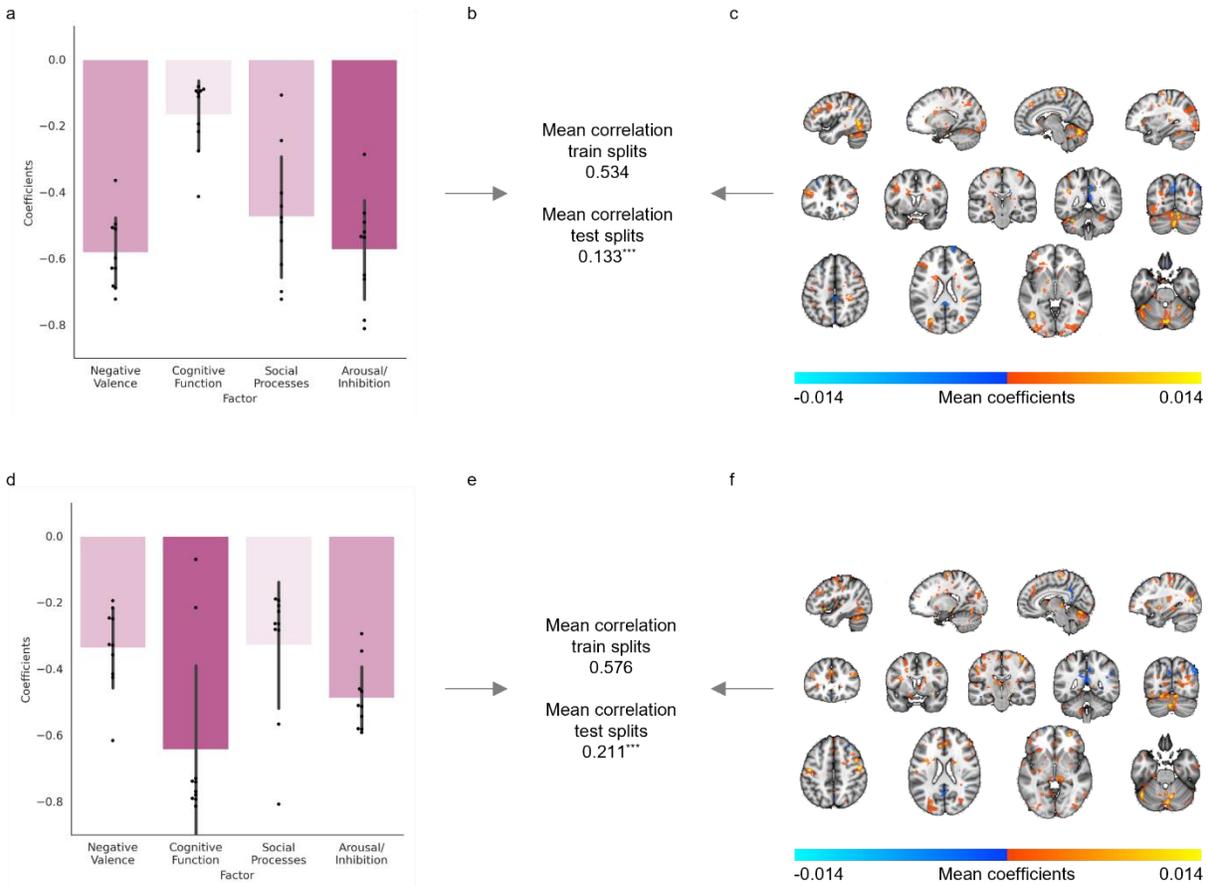
256 We then aimed to determine whether multivariate patterns of deviation from the reference models were
257 associated with cross-diagnostic symptomatology. To achieve this, we input whole-brain (unthresholded)
258 deviation maps and factor loadings for negative valence, cognitive function, social processes and
259 arousal/inhibition domains from prior work¹⁷ to an established penalised canonical correlation analysis
260 (CCA) framework that enforces sparsity (sparse CCA, SCCA; functional domain loading scores were
261 available for 217 patients)^{18,19}. Significant out of sample associations (10 fold 70% - 30% training- test split)
262 were detected both for faces>shapes and faces>baseline contrasts (mean r of test splits 0.133 and 0.211
263 respectively, both $p < 0.001$ by permutation test; Fig. 6b, e) but with distinct patterns of effects both in terms
264 of symptom domains and associated brain regions. More specifically, for the faces>shapes contrast,
265 decreased functioning predominantly in the negative valence and arousal/inhibition domains (Fig. 6a) was
266 associated with a pattern of deviations including the right insula, the bilateral medial prefrontal cortex and
267 pre- and post- central gyri, the bilateral inferior temporal gyrus, lingual gyrus, bilateral hippocampus and
268 the right thalamus, as well as the regions in the medial and left lateral cerebellum (Fig. 6c). By comparison,
269 for the faces>baseline contrast factor loadings for cognitive functioning and arousal/inhibition (Fig. 6d) were
270 most strongly related to a pattern comprising bilateral insula, the anterior-to-medial cingulate cortex
271 extending to the dorsal medial prefrontal cortex, the pre- and post- central gyri, the right middle frontal and
272 bilateral inferior frontal gyrus, and the bilateral hippocampus, caudate, putamen and amygdala, and the

273 medial and left lateral cerebellum (Fig. 6f). The SCCA was repeated to relate participant's diagnoses with
274 their whole-brain (unthresholded) deviation maps. In contrast to the cross-diagnostic symptom domains,
275 there was no association between diagnostic labels and deviation scores. Mean canonical correlations were
276 small (mean r of test splits <0.09 for both faces>shapes and faces>baseline models), and this was not
277 statistically significant as determined by 1000-fold permutation testing.
278
279



280

281 **Figure 5: Testing the faces>shapes (left) and faces>baseline normative models with the MIND-Set cohort.** Normative
282 Probability Maps illustrate the percentage of participants of the clinical sample who had positive (hot colours) or negative deviations
283 (cool colours) $> \pm 2.6$ within each voxel, for the faces>shapes (a) and faces>baseline (d) models. Histograms and box plots show the
284 relative frequency and mean number of the total deviations that a participant has for faces>shapes (b,c) and faces>baseline (e,f)
285 models. Positive: x,y,z = 9,-16,-16, Negative: x,y,z = 5,-29,-24.



286

287 **Figure 6: Sparse canonical correlation analyses (SCCA) between functional domains, and deviation scores from**
 288 **faces>shapes or faces>baseline normative models.** Weights per factor to latent variable of psycho-social functioning (a,d).
 289 Canonical correlation between 4 functional domains and whole-brain deviation scores from (b) faces>shapes and (e) faces>baseline
 290 normative models (regularisation 10%). Mean voxel-wise weights to latent variable of deviation scores from (c) faces>shapes
 291 normative models and from (f) faces>baseline. All results are statistically significant with 1000-fold permutation tests (*** = $p < 0.001$).
 292 $x, y, z, = [-42, -17, 8, 33], [29, 4, -21, -46, -71], [47, 22, -3, -28]$.
 293
 294

295

296

296 **Spatial extent of deviations highlights similarities across, and differences between diagnoses**

297 Finally, we were interested in mapping the spatial distribution of the deviations within the clinical sample,
 298 and whether this varied according to the participant's mental health diagnosis or neurodivergence (note
 299 that subjects can be in multiple categories; Sup Fig. 5). For each diagnosis, the pattern of deviations was
 300 highly heterogeneous, providing further support for high degree of inter-individual heterogeneity we have
 301 reported previously for mental disorders^{4,5,7}, and underlining the need to move beyond case-control
 302 comparisons at the level of diagnostic groups.

303

304

305 **Discussion:**

306 In this study we made use of six large publicly available datasets of participants completing the fMRI EFMT
307 to build a reference normative model of functional activation underlying emotional face processing. We
308 collated data from over 7500 participants and show that our voxel-wise models can explain up to 50% of
309 variance in observed BOLD signal, with the remaining unexplained variance representative of individual
310 differences in functional activation (deviation scores). We unpacked the variance explained by the models,
311 to show how the predicted activation related to the models' input variables, namely demographics,
312 variations in task design, and acquisition parameters. Lastly, we tested our reference model with data from
313 a sample of patients with heterogenous and frequently co-occurring psychiatric conditions (mood and
314 anxiety disorders, and neurodevelopmental conditions). Our analyses show that: (i) there is considerable
315 inter-individual variation superimposed on the group effects customarily reported in fMRI studies, (ii) that
316 such variation is predictive of psychiatric symptom domains in a cross-diagnostic fashion and (iii) while an
317 overall effect of diagnosis was evident, this was highly individualised in that the overlap of deviations
318 amongst individuals with the same diagnosis was low. This implies that there are brain regions wherein
319 patients more frequently have deviations irrespective of the type of diagnoses, and other regions wherein
320 the frequency of deviations appears specific to the mental-health condition or neurodivergent diagnosis.

321
322 A key feature of the normative modelling framework in the context of multi-site fMRI data is that it allows us
323 to aggregate data across multiple samples by binding them to a common reference model. This provides
324 multiple benefits: it removes site effects from the data without requiring the data to be harmonized²⁰, which
325 avoids the introduction of certain biases due to harmonisation²¹ and allows meaningful comparisons to be
326 drawn across studies. For example, this allows aggregation of different studies to better understand
327 variation across cohorts or across the lifespan and to understand the effect of different task parameters on
328 functional activity across cohorts. Moreover, by placing each individual within the same reference model
329 this provides the ability to quantify, compare and ultimately parse heterogeneity across studies.

330
331 Traditional group-level task contrasts, as shown in Fig. 1, inform us of the region's most consistently
332 activated across participants/groups during task conditions. Their interpretation has relied heavily on the
333 assumption of spatial homogeneity of activation between subjects; an assumption that the deviation scores
334 from our reference model show to be largely untrue (Fig. 2). We show that such group effects reflect a small
335 proportion of the variation amongst individuals and using the normative modelling framework we map the
336 underlying heterogeneity, separating variation in the intensity and spatial extent of task-evoked functional
337 activation between-subjects attributable to known factors such as site effects, demographics, acquisition
338 parameters, and differences in EFMT paradigm design. More importantly, we show that residual differences
339 in the neuronal effects elicited by the task are highly meaningful in that they were predictive of psychiatric

340 symptomatology and can be used to understand inter-individual differences in functional anatomy and its
341 relation to clinical variables. In our test reference population, while every voxel of the brain had at least one
342 participant with a large deviation, some regions considered active during the faces condition (as compared
343 to shapes), such as the medial occipital lobe, fusiform gyrus and inferior temporal lobe, were also regions
344 in which positive deviations were frequently observed.

345 When building our reference models, we chose to include and control for multiple variables that we
346 reasoned may influence the BOLD signal observed. These included demographic factors such as age and
347 sex, and task design choices that could influence the BOLD signal generated, as well as acquisition
348 parameters that could influence the BOLD signal recorded. Some effects, such as that of age and task
349 instructions were relatively strong and interpretable, for example, increased age predicted decreased
350 activity in surface areas of the brain and regions surrounding the ventricles likely reflecting decreased signal
351 due to age-related atrophy, and instructing participants to match emotional expressions, as opposed to
352 matching faces, increased the predicted activity in the thalamus which may reflect increased engagement
353 of regions associated with affective processing. On the other hand, other variables explained relatively little
354 variance in the predictions (e.g. sex). In our sample many predictor variables were collinear across sites
355 which limited our ability to detect systematic differences resulting for example from differences in the task
356 paradigm. In this work we decided to keep all variables in the model and used structure coefficients to
357 identify the importance of different variables, which are relatively insensitive to collinearity. This follows prior
358 work to identify specific effects of input variables on model predictions, for example the influence of specific
359 adversity types on predicted morphometric changes (Holz et al., *in prep*). However future researchers may
360 consider reducing the dimensionality of their inputs prior to model construction. Future studies with larger
361 numbers of more diverse samples (e.g. more variations on the basic task design), as is possible in
362 consortium such as ENIGMA, will allow for more fine grained analyses of the effect of task parameters on
363 inter-individual variation within the population.

364

365 We demonstrated that distinct patterns of deviations, derived from each model type (faces>shapes or
366 faces>baseline), were associated with unique profiles of functioning across four transdiagnostic domains.
367 The distinct patterns of effects, in terms of the implicated symptom domains and associated brain regions,
368 make sense in the context of relevant existing literature. For example, negative affect, impulsivity and
369 emotional liability have previously been related to functional activity within the bilateral insula, motor cortex
370 and hippocampus²², and cognitive functioning has been linked to activity within the medial prefrontal cortex,
371 anterior-to-medial cingulate cortex, superior frontal gyrus. This not only validates the interpretability of
372 findings from these normative modelling analyse, but also illustrates the potential for future researchers to
373 use individualised deviation maps to better understand the neural processes that underly cognitive and
374 affective functioning, within and across diagnostic boundaries. Furthermore, approaching dysfunction
375 through the normative modelling framework and transdiagnostic functional domains appears to more

376 closely relate to underlying biology. This reflects practitioners implementation of clinical care and the use
377 of overlapping treatments for differing disorders, which often does not fit a binary classification paradigm.
378 Using this modelling approach may also better allow for the quantification of neurodivergence, not as being
379 'disordered' but rather as varying phenotypic expressions along a characterised spectrum.

380 It should be noted, however, that within any one voxel of the brain, only ~20% of the clinical sample
381 (be that in the total sample, or within disorders) had large deviations. This suggests that the exact location
382 of deviations is very variable between individuals, and could explain why many prior studies have not found
383 significant differences when performing traditional case-control analyses. In this study, we aimed to
384 estimate the degree to which the deviations from the normative models were associated with cross-
385 diagnostic symptomatology, but other approaches may also be useful, as outlined in Rutherford, et al.²³.
386 For example, clustering algorithms could be applied to derive a stratification of individuals^{4,5} or to identify
387 heterogenous yet convergent functional processes (many-to-one functional mappings)²⁴, and supervised
388 learning methods may be useful to assess the degree to which specific clinical variables can be predicted
389 from the patterns of deviation we report.

390
391 Interestingly, the normative models of functional activation built using the faces>shapes have a different
392 pattern of association with symptomatology relative to the faces>baseline contrast. This suggests that the
393 two contrasts carry complementary information about psychopathology. The frequency of deviation scores
394 was significantly greater in the clinical cohort, compared to the reference cohort, when using the
395 faces>shapes contrast, and the weights attributed to each of these deviations (at a voxel-level) in the SCCA
396 were associated with different symptom domains. Neither contrast was significantly predictive of diagnosis.
397 By comparison, the relationship between the frequency of deviation scores and domains of function was
398 stronger when using models built using the faces>baseline contrast, which was further supported by the
399 stronger canonical correlation between factor loadings for functional domains and deviations from the
400 faces>baseline models. Taken together, this could be interpreted to suggest that widespread deviations,
401 best captured by the faces>shapes contrast, are indicative of global alterations in functioning which are
402 broadly linked to different clinical diagnoses. By comparison, fewer but more focal deviations and/or the
403 ability to detect abnormal baselines of activation^{25,26}, best revealed using the faces>baseline, have greater
404 relation to specific functional domains. Future researchers should carefully consider the task contrast used
405 to construct their normative models.

406
407 Concerns for the within-subject reliability of task-based fMRI data²⁷ are not to be dismissed in the context
408 of our models which are currently built on cross-sectional data. While we acknowledge the limitations
409 imposed due to the limits of test-retest reliability of task fMRI, our results nevertheless provide encouraging
410 evidence for the use of task fMRI readouts as individualised biomarkers as we show by their ability to predict

411 clinical variables in the context of SCCA. While we were not able to assess test-retest reliability directly
412 within the context of normative models due to a lack of test-retest data for the EFMT, one great strength of
413 the normative modelling approach is the shift from looking only at mean activity, to estimating the underlying
414 variance. This can implicitly down-weight regions or individuals that are less reliable. Indeed, it could be
415 expected that repeated sampling would fall within the same variance distribution, and as such our broad
416 understanding of brain function remains quite stable. The normative modelling framework is also ideally
417 positioned to directly test the reproducibility of fMRI within subjects. In follow-on work to the present
418 manuscript, we are currently developing an extension to explicitly include test-retest variability in the model
419 by testing reference models with repeat scans from participants, and compare individuals' deviation scores
420 between the two tests, whilst explicitly quantifying within subject variance, such that it provides a lower
421 bound on the size of deviation that can be considered meaningful (Bučková et. al., *in prep*). Alternatively,
422 where multiple repeats are available, hierarchical models can be used to accommodate dependencies
423 between subjects²⁰ which would provide more precise estimates of individual deviations. The application of
424 the normative modelling method to fMRI can easily be generalised to other tasks (e.g. the monetary
425 incentive delay incentive processing task or n-back work memory task) and need not stop at predicting
426 functional activation. With the right data sets, this method could use fMRI data to predict many other
427 variables including psychophysiological responses or subjective ratings of affect.

428

429 With this work, we show the potential for the normative modelling framework to be applied to large task-
430 based fMRI data sets to bind heterogeneous datasets to a common reference model and enable meaningful
431 comparisons between them. Using this approach, we illustrate the heterogeneous intensity and spatial
432 location of task-evoked activation within the general population² using the EFMT in a sample of over 7500
433 participants. Further, we applied this model to patients with a current diagnosis (mood and anxiety
434 disorders, ASD and/or ADHD) and demonstrate the transdiagnostic clinical relevance and further potential
435 for deviation scores derived from this method. The potential of this method is clear; normative modelling of
436 task-based functional activation can facilitate a better understanding of individual differences in complex
437 brain-behaviour relationships, and further our understanding of how these differences relate to mental
438 health and neurodivergence.

439

440 **Methods:**

441 **Data sets:** We collated a large reference sample from 6 independent sites for whom high quality fMRI
442 data for the EFMT are available: AOMIC PIOP2, Duke Neurogenetics Study, HCP Development, HCP
443 Young Adult (1200 release), UK Biobank, and the MIND-Set cohort which also includes a clinical
444 population. For sample details per site see Table 1.

445
446 **fMRI task paradigms:** All sites collected a variant of the EFMT⁹. Although specific parameters varied, the
447 overall design was consistent: in each face trial participants were presented with three images of human
448 faces in a triangular formation. Participants were instructed to identify which of two faces/expressions
449 presented at the bottom of the screen matched the one presented at the top of the screen by pushing a
450 button with the index finger of their left or right hand. Multiple face trials were presented in one face block,
451 and the task included multiple face blocks (see Table 1 for the number of trials per block, and number of
452 blocks per site). As a somatomotor control, participants also completed shape trials, wherein they were
453 presented with three geometric shapes (circles and ovals) and asked to indicate which of the two shapes
454 presented at the bottom of the screen matched the one at the top. Multiple shape trials were concatenated
455 to form one shape block, which were interspersed between face blocks.

456 Two paradigms (HCP Young Adult and HCP Development) included an inter-trial interval (white
457 fixation cross on black screen), and three sites (HCP Young Adult, HCP Development and AOMIC PIOP2)
458 had an instruction trial that preceded the start of each block. Tasks varied in their duration from 150 to 290
459 seconds, which indirectly corresponded to the acquisition of between 135 and 336 functional volumes.

460
461 **fMRI data acquisition:**

462 Site specific acquisition parameters per site are detailed in Table 1, and in the following site specific
463 protocols: AOMIC PIOP2¹⁵, HCP Young Adult¹³, HCP Development²⁸, UKBiobank²⁹, Duke Neurogenetics
464 Study (<https://www.haririlab.com/methods/amygdala.html>) and MIND-Set³⁰.

465
466 **fMRI pre-processing:** Data pre-processing was harmonised across all sites; a FSL-based pipeline³¹ was
467 consistently applied to decrease the likelihood of introducing variance due to pre-processing differences.
468 Since the HCP young adult, HCP development and UKB Biobank data were already processed relatively
469 consistently, we reused the processing pipelines provided by the respective consortia (for HCP sites we
470 used the minimal processing pipeline)^{29,32}, with additional steps taken as necessary (e.g. matching
471 smoothing kernels across studies). At a within-subject level, all functional data underwent gradient
472 unwarping, motion correction, fieldmap-based EPI distortion correction (where fieldmaps were available),
473 boundary-based registration of EPI to structural T1-weighted scan, denoising for secondary head motion-
474 related artifacts using automatic noise selection, as implemented in ICA-AROMA³³, non-linear registration
475 into MNI152 space, and grand-mean intensity normalization. All data were spatially smoothed using a 5
476 mm FWHM Gaussian kernel.

477
478 **Quality control:** Participants were excluded if their mean relative RMS was greater than 0.5mm. Additional
479 quality control was performed for signal coverage in the prefrontal cortex for the UK Biobank sample (see
480 supplementary methods).

481
482 **fMRI general linear modelling (GLM) – single subject:** We matched the methodological approach used
483 to estimate the parameters within a GLM-based analysis, given evidence to suggest this analytic step can
484 significantly contribute to the variability of reported results between sites³⁴. Therefore, for each site, the
485 linear model parameter were estimated using the FSL software package version 6.03 (HCP Young Adult,
486 HCP Development, MIND-Set, Duke Neurogenetics Study; <http://fsl.fmrib.ox.ac.uk/>) or as downloadable
487 form UK Biobank²⁹. Two regressors were constructed from the faces and shapes blocks which were then
488 convolved with a canonical double-gamma haemodynamic response function and combined with the
489 temporal derivatives of each main regressor. These were treated as nuisance regressors and served to
490 accommodate slight variations in slice timing or in the haemodynamic response. Data were pre-whitened
491 using a version of FSL-FILM customized to accommodate surface data, the model and data were high-
492 pass filtered (200s cut-off). Fixed-effects GLMs were estimated using FSL-FLAME 1: first for independent
493 runs, then when necessary combining two runs into a single model for each participant (HCP Young Adult).
494 and the AOMIC, DNS and MIND-Set maps were transformed into standard space using FNIRT³⁵. We
495 created summary group level maps per site (for a random sample of 100 participants), as a sanity check to
496 ensure the data was otherwise comparable to past literature, and performed a case-control comparison
497 between patients with a current diagnosis (mood and anxiety disorders, ASD and/or ADHD) and unaffected
498 controls in the MIND-Set cohort.

499
500 **Normative models:** The z-statistic maps from the contrast face>shapes (5mm smoothed in standard
501 space), for each subject, were used as response variables for the normative models. That is, we specified
502 a functional relationship between a vector of covariates and responses. We created normative models of
503 EFMT-related BOLD activation maps, as a function of sex, age, acquisition and task parameters (task
504 duration (s), number of target blocks, instructions given to participants, the task stimuli), by training a
505 Bayesian Linear Regression (BLR) model to predict BOLD signal for the faces>shapes contrast.
506 Generalisability was assessed by using a half-split train-test sample (train: n = 3877, test: n = 3764). In
507 preliminary analyses, we compared a warped model which can model non-Gaussianity with a vanilla
508 Gaussian BLR model. Since the fit was comparable across most metrics and regions, we focus on the
509 simpler Gaussian model below. We included dummy coded site-related variables as additional covariates
510 of no-interest. We also created models to predict BOLD signal for the faces condition alone (i.e.
511 face>baseline contrast). This was performed in the Predictive Clinical Neuroscience toolkit (PCNtoolkit)
512 software v0.26 (<https://pcntoolkit.readthedocs.io/en/latest>) implemented in python 3.8.

513

514 **Quantifying voxel-wise deviations from the reference normative model:** To estimate a pattern of
515 regional deviations from typical brain function for each participant, we derived a normative probability map
516 (NPM) that quantifies the voxel-wise deviation from the normative model. The subject-specific Z-score
517 indicates the difference between the predicted activation and true activation scaled by the prediction
518 variance. We thresholded participant's NPM at $Z = \pm 2.6$ (i.e. $p < .005$)⁷ using `fslmaths` and summed the
519 number of significantly deviating voxels for each participant, and then across the total sample.

520
521 **Effects of input variables on model predictions:** In order to probe the magnitude of the influence of task
522 design parameters on the predicted BOLD signal, we illustrated the structure coefficients (correlation
523 coefficients) of each task parameter-related variable (task duration (s), number of target blocks, instructions
524 given to participants, the task stimuli), as well as for age, sex and ICV. This approach is preferable to
525 regression coefficients when variables are collinear³⁶.

526
527 **Clinical application:** We tested the normative models we created using the reference data, with a
528 heterogeneous patient sample from the MIND-Set cohort ($n = 236$, mean age = 37.1 ± 13.27 ; 41.94%
529 female). This is a naturalistic and highly co-morbid sample derived from out-patients of the psychiatry
530 department of Radboud University Medical Centre. This included 150 patients diagnosed with a current
531 mood disorder (unipolar or bipolar depressive disorder), 12 with generalised anxiety disorder, 22 with social
532 phobia, 14 with panic disorder, 71 with attention-deficit-hyperactive-disorder, and 55 autistic individuals
533 (see Table 1 for full details and note that subjects can be in multiple diagnostic categories). The clinical
534 relevance of our models can also be tested in the context of transdiagnostic symptom domains; a
535 conceptualisation of mental functioning that transcends diagnostic boundaries and allows for nuanced
536 brain-behaviour interpretations. As such, for 217 (of our 236) patients for whom all required data was
537 available, we repeated a previously validated factor analysis method (performed in SPSS v24.0, oblique
538 rotation)¹⁷ to obtain individual factor loadings on 4 functional domains: (1) negative valence, (2) cognitive
539 function, (3) social processes and (4) arousal/inhibition.

540
541 **Quantifying patients' voxel-wise deviations from the reference normative model:** As for the reference
542 cohort, we generated NPMs to estimate the pattern of regional deviations from typical brain function for
543 each participant, and summed across the sample. We then used a Mann-Whitney U test to compare the
544 frequency of deviations ($> \pm 2.6$) between the reference controls and patients from the MIND-Set cohort.

545
546 **Relating deviations to transdiagnostic functional domains:** In order to map the association of the
547 deviation scores with cross-diagnostic symptomatology, we performed sparse canonical correlation
548 analyses (SCCA) to relate participant's scores in the four aforementioned functional domains or their
549 diagnoses, to their whole-brain (unthresholded) deviation maps using an established penalised CCA
550 framework that enforces sparsity^{18,19}. Specifically, we applied variable shrinkage by adding an l_1 -norm

551 penalty term to stabilise the CCA estimation and ensure the weights for the deviation scores were more
552 interpretable. We follow the formulation outlined in Witten, et al. ¹⁸, where we refer to for details. In brief,
553 given two data matrices X and Y with dimensions $n \times p$ and $n \times q$ respectively (here, these are the cross-
554 diagnostic factor loadings and whole-brain deviations), and two weight vectors \mathbf{u} and \mathbf{v} this involves
555 maximising the quantity $\rho = \mathbf{u}^T X^T Y \mathbf{v}$ subject to the constraints $\|\mathbf{u}\|_2^2 \leq 1$, $\|\mathbf{v}\|_2^2 \leq 1$, $\|\mathbf{u}\|_1 \leq c_1$ and $\|\mathbf{v}\|_1 \leq$
556 c_2 , where the penalties $p(\mathbf{u})$ and $p(\mathbf{v})$ are the standard L1-norm. We set the regularisation parameters for
557 each view heuristically ($c_1 = 0.9p$ corresponding to light regularisation for the factor scores, $c_1 = 0.1q$,
558 corresponding to heavy regularisation for the deviation maps such that no more than 10% of voxels were
559 selected). While it is possible that better performance would be obtained by optimising the regularisation
560 parameters across a grid, we did not pursue that here due to the moderate sample size for the clinical
561 dataset. We assessed generalisability of SCCA by splitting the data in to 70% training data and 30% test
562 10 times. Finally, we wrapped the entire procedure in a permutation test where we randomly permuted the
563 rows of one of the matrices 1000 times to compute an empirical null distribution for significance testing.

564

565 ***Spatial patterns of deviations by primary and co-occurring diagnoses:*** We illustrated the spatial
566 heterogeneity in deviations between different diagnoses (note that subjects can be in multiple categories),
567 and further, compared patients with a single diagnoses to those with two, three, or more than three
568 diagnoses, to determine whether and if so, how the location of deviations related to the number of co-
569 occurring diagnoses a patient has.

570

571 **Data availability:**

572 Scripts for running the analysis and visualizations are available on GitHub ([https://github.com/predictive-](https://github.com/predictive-clinical-neuroscience/EFMT_Norm_Models)
573 [clinical-neuroscience/EFMT_Norm_Models](https://github.com/predictive-clinical-neuroscience/EFMT_Norm_Models)).

574

575 **Acknowledgements:**

576 The Duke Neurogenetics Study was supported by Duke University and National Institutes of Health (NIH)
577 grants R01DA031579 and R01DA033369. Ahmad R. Hariri (Duke University) is further supported by NIH
578 grant R01AG049789. The Duke Brain Imaging and Analysis Center's computing cluster, upon which all
579 DNS analyses heavily rely, was supported by the Office of the Director, NIH, under Award Number
580 S10OD021480. CFB gratefully acknowledges funding from the Wellcome Trust Collaborative Award in
581 Science 215573/Z/19/Z and the Netherlands Organization for Scientific Research Vici Grant No. 17854 and
582 NWO-CAS Grant No. 012-200-013.

Table 1: Sample details, functional scan acquisition parameters and Emotional Face Matching Task parameters for data included in the normative models.

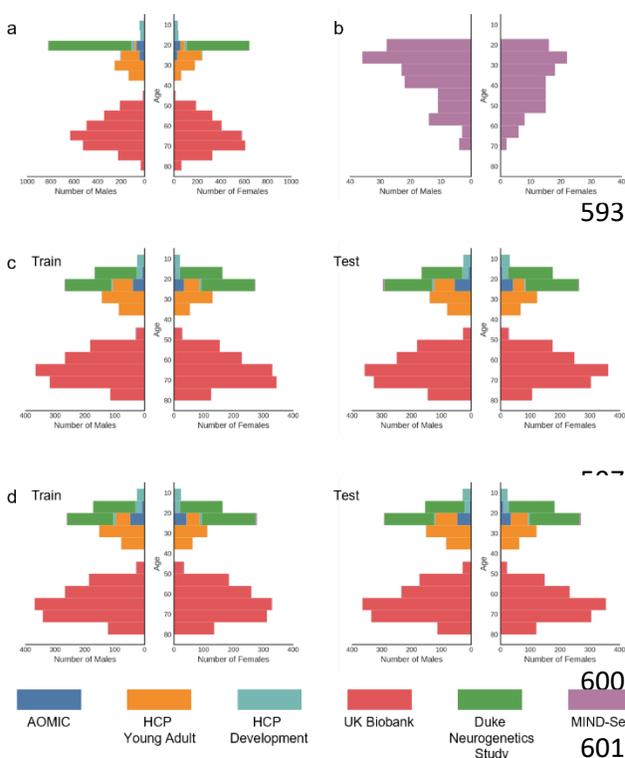
Sample details				Functional scan acquisition parameters				Emotional Face Matching Task parameters									
Site	Sample size	Sex	Age (mean+ S.D) [range]	Scanner	TE/TR (ms)	Multi-band Factor	Flip angle	Matching Rule /Instructions	Target Stimulus	Trials per block/ Blocks/ Total Trials	Trial duration (s)	Instruction duration (s)	Inter-trial interval (s)	Block duration (s)	Task duration (s)	Inter-block interval (s)	Volumes acquired
Human Connectome Project Young Adult	1044	561 F (53.7%)	28.76±3.70 [22-37]	3T Siemens Skyra	33.1/720	8	52	Match faces: Decide which of two faces presented on the bottom of the screen match the face at the top of the screen.	Angry and fearful faces: Nim-Stim Face Stimulus Set	6/3/18	2	3	1	21	156	NA	176
Human Connectome Project Development	201	110 F (54.7%)	13.86±3.83 [8-21]	3T Siemens Prisma	37/800												178
UK Biobank	5000	2487 F (49.7%)	63.99±7.45 [46- 82]	3T Siemens Skyra	39/735					Match faces: Indicate which face [or shape] on the bottom row matches the face on the top row.	NA/5/NA	NA	NA	NA	253	8	366
Amsterdam Open MRI Collection Population Imaging of Psychology	200	114 F (57.0%)	22.16±1.79 [18.25– 26.25]	3T Phillips Achieva dStream	28/2000	NA	76.1	Match expression: Match the emotional expression of the target face as quickly as possible.	Angry, fearful, surprised, neutral faces: Ekman and Friesen, 1976	6/4/24	when selected or up to 4.8s	10	5s – Reaction Time	~25	290	5	135
Duke Neurogenetics Study	1246	707 F (56.7%)	20.22±1.21 [18.09-23.07]	3T GE MR750	30/2000												NA
MIND-Set	Reference: 37/309 Clinical Test: 36/309	21 F (56.7%) 21 F (58.3%)	38.0±16.11 [20-74] 37.1±16.50 [20-70]	3T Siemens Magentom Prisma	34/1000	6	60	Match expression: Indicate which one of the bottom two faces matched the top face in terms of emotional expression.	Angry and fearful faces: Nim-Stim Face Stimulus Set	6/2/12	5	NA	NA	30	150	NA	150
Additional details of clinical samples:	Sample size	Sex	Age (mean+ S.D) [range]	Current diagnoses					Number of Diagnoses (% of total sample)								
MIND-Set	236/309	99 F (41.9%)	37.1±13.27 [20-74]	150 Mood Disorder 71 Attention deficit hyperactivity disorder (ADHD) 55 Autism spectrum disorder (ASD) 22 Social Phobia 14 Panic Disorder	12 Generalised Anxiety Disorder 7 Anxiety disorder NOS* 6 Obsessive Compulsive Disorder 5 Post Traumatic Stress Disorder 4 Specific Phobia 2 Agoraphobia	1: 66 (27.9%) 2: 65 (27.5%) 3: 39 (16.5%) >3: 8 (3.38%)											

Note: Underline indicates that this parameter was input as a variable in the normative models. *NOS (Not otherwise specified).

585

SUPPLEMENTARY MATERIALS

586 SUPPLEMENTARY METHODS – *Sample Details:*



Supplementary Figure 1: Age and sex distributions of (a) the total reference sample, (b) the total clinical test sample (MIND-Set) (c) the faces>shapes train (left) and test (right) split, and the (d) the faces>baseline train (left) and test (right) split.

602 SUPPLEMENTARY METHODS – *Signal coverage of the prefrontal cortex (PFC):*

603 Due to air-tissue inhomogeneities which can diminish the acquired BOLD signal to such a degree that no
604 activations are visible, a notorious effect within the ventral PFC, we performed targeted quality control for
605 this extended region. Binary ROI masks were created for the dorsal ventro-medial PFC (d-vmPFC), ventral
606 vmPFC (v-vmPFC), lateral vmPFC (l-vmPFC) and the dorso-medial PFC (dmPFC), as defined by the AAL2
607 atlas regions 25, 27, 33 and 1 respectively (see Supplementary Figure 1C). The percentage of voxels with
608 an absolute value greater than 0 for the contrast faces > shapes within each ROI was determined (i.e.
609 where any signal was present regardless of its relative direction; see Supplementary Figure 1A,B). While
610 most sites had good coverage, the coverage within the ventral and lateral vmPFC regions were particularly
611 variable for the UK Biobank data. We therefore performed this step *only on data from the UK Biobank site*;
612 this selectivity was made possible by the large number of participants we had access to, and our need to
613 include but a fraction of the total available sample. We ranked participants in descending order of the
614 percent of their v-vmPFC, l-vmPFC, d-vmPFC, and dmPFC covered, respectively, and selected the first
615 5000 participants. We also collected the percentage covered value for a bilateral amygdala ROI mask, but
616 made no exclusion/inclusions on this basis as coverage was very high across all participants and all sites.
617

618 **Supplementary Table 1: Motion QC**

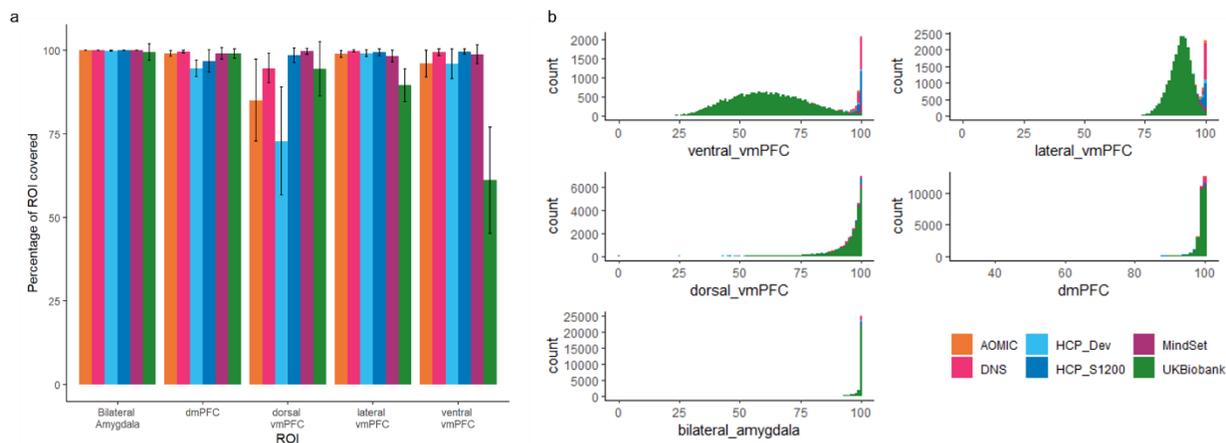
Site	Full sample size	Included	Excluded
AOMIC	217	217	0
DNS	1263	1246	17
HCP Young Adult	1044	1044	0
HCP Development	256	256	0
UK Biobank	26167	5000	N/A
MIND-Set	393	389	4

619

620 **Supplementary Table 2: vmPFC QC**

ROI	Site	Sample size (n)	Mean percentage of ROI covered	Standard deviation
Bilateral Amygdala	AOMIC	217	100	0
	DNS	1246	100	0
	HCP Development	256	99.94936	0.162845
	HCP Young Adult	1044	99.98188	0.098589
	MIND Set	389	99.99966	0.006707
	UK Biobank	26120	99.47363	2.399575
dmPFC	AOMIC	217	99.10345	0.830591
	DNS	1246	99.61285	0.411399
	HCP Development	256	94.6524	2.460686
	HCP Young Adult	1044	96.85686	3.318827
	MIND Set	389	99.069	1.785034
	UK Biobank	26120	99.0493	1.405906
Dorsal vmPFC	AOMIC	217	85.08059	12.23881
	DNS	1246	94.68589	4.406996
	HCP Development	256	72.90154	16.11571
	HCP Young Adult	1044	98.58104	2.192929
	MIND Set	389	99.75218	0.886996
	UK Biobank	26120	94.50698	8.107753
Lateral vmPFC	AOMIC	217	98.90096	1.063159
	DNS	1246	99.82392	0.33135
	HCP Development	256	99.14376	1.003451
	HCP Young Adult	1044	99.46709	0.978404
	MIND Set	389	98.23297	1.745961
	UK Biobank	26120	89.56122	4.919904
Ventral vmPFC	AOMIC	217	96.06479	3.960334
	DNS	1246	99.45099	0.964661
	HCP Development	256	95.999	4.42593
	HCP Young Adult	1044	99.63665	0.85699
	MIND Set	389	98.83508	2.868568
	UK Biobank	26120	61.17452	15.9091

621



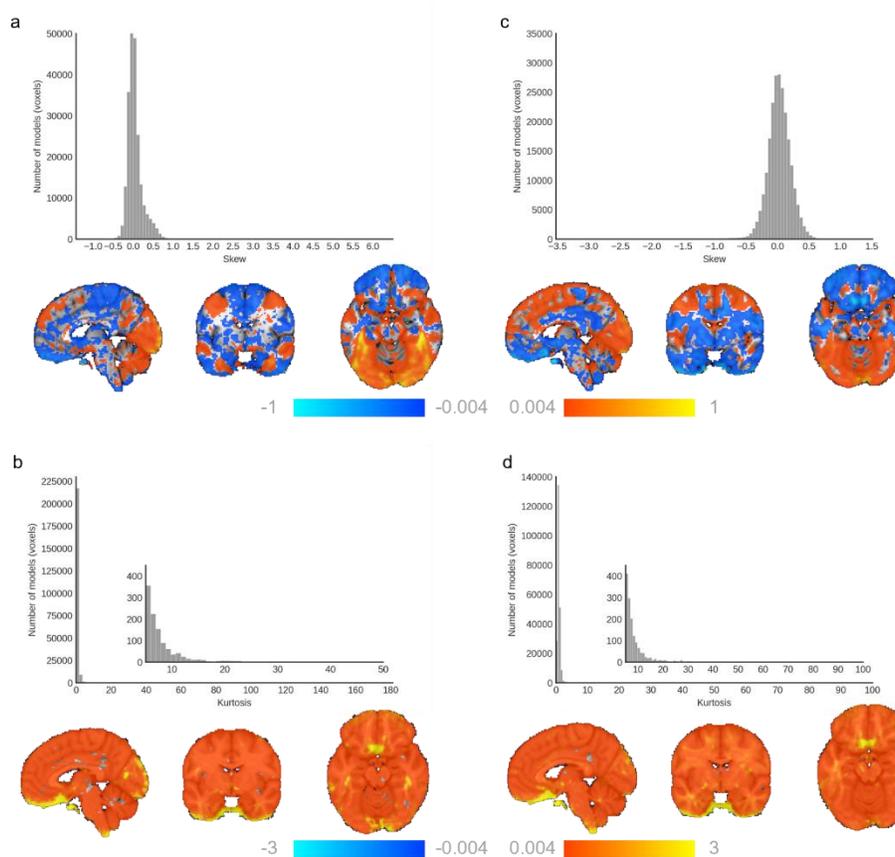
622

623 **Supplementary Figure 2: vmPFC QC metrics.** (a) Mean percentage of each ROI with signal greater than 0, used for quality control.
 624 Error bars show +/- standard deviation (b) Stacked histograms (raw participant count) of the percentage of each ROI covered, coloured
 625 by site.

626

627 **SUPPLEMENTARY RESULTS – Evaluation of reference normative models:**

628

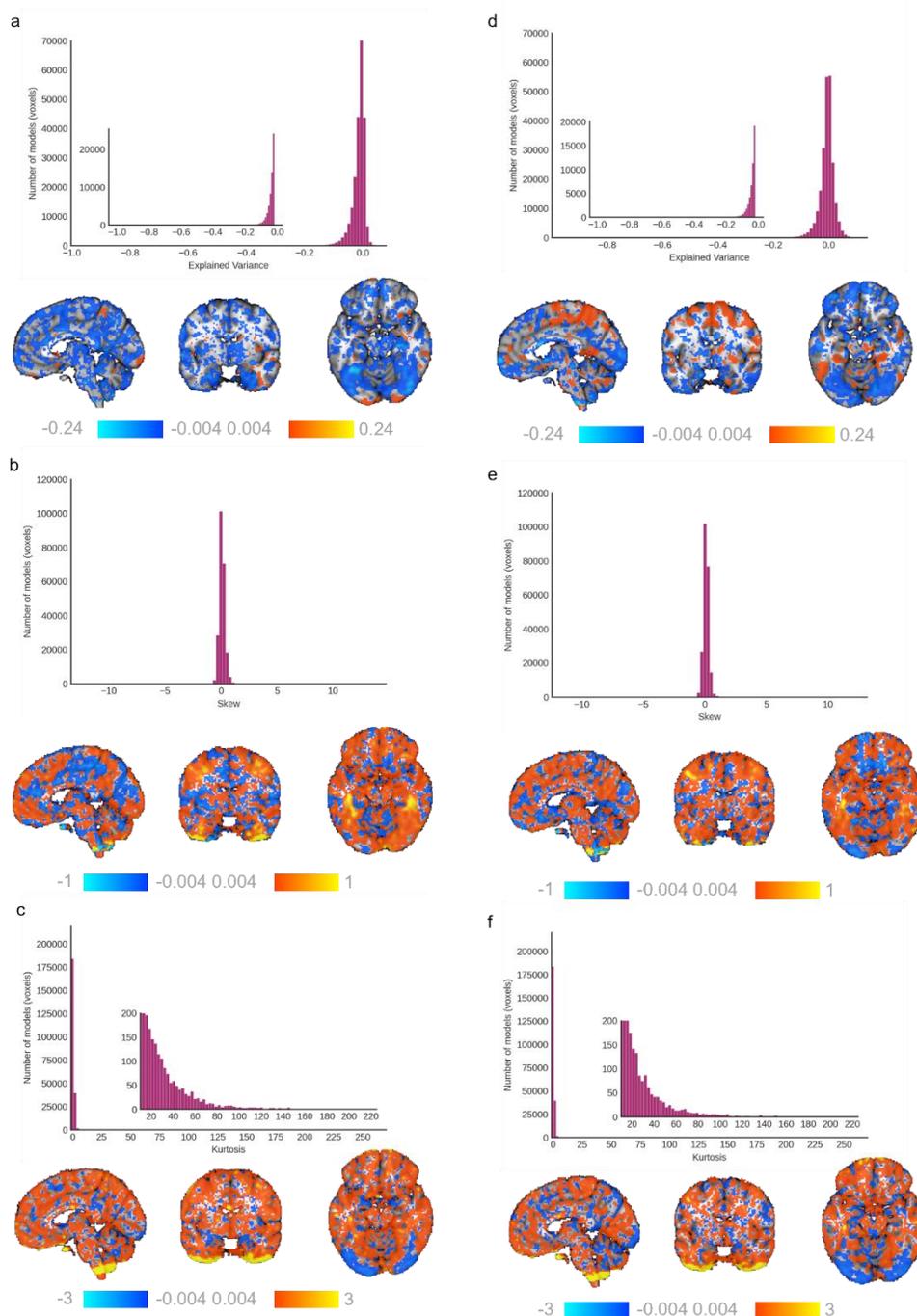


629

630 **Supplementary Figure 3: Evaluation of the faces>shapes (left) and faces>baseline (right) reference normative models.**
 631 Histograms show the skew (a,c), and kurtosis (b,d) of the normative models, and their respective illustration on the brain (x,y,z= 4,-6,-
 632 15).
 633

634 **SUPPLEMENTARY RESULTS – Evaluation of normative models when applied to MIND-Set cohort:**

635



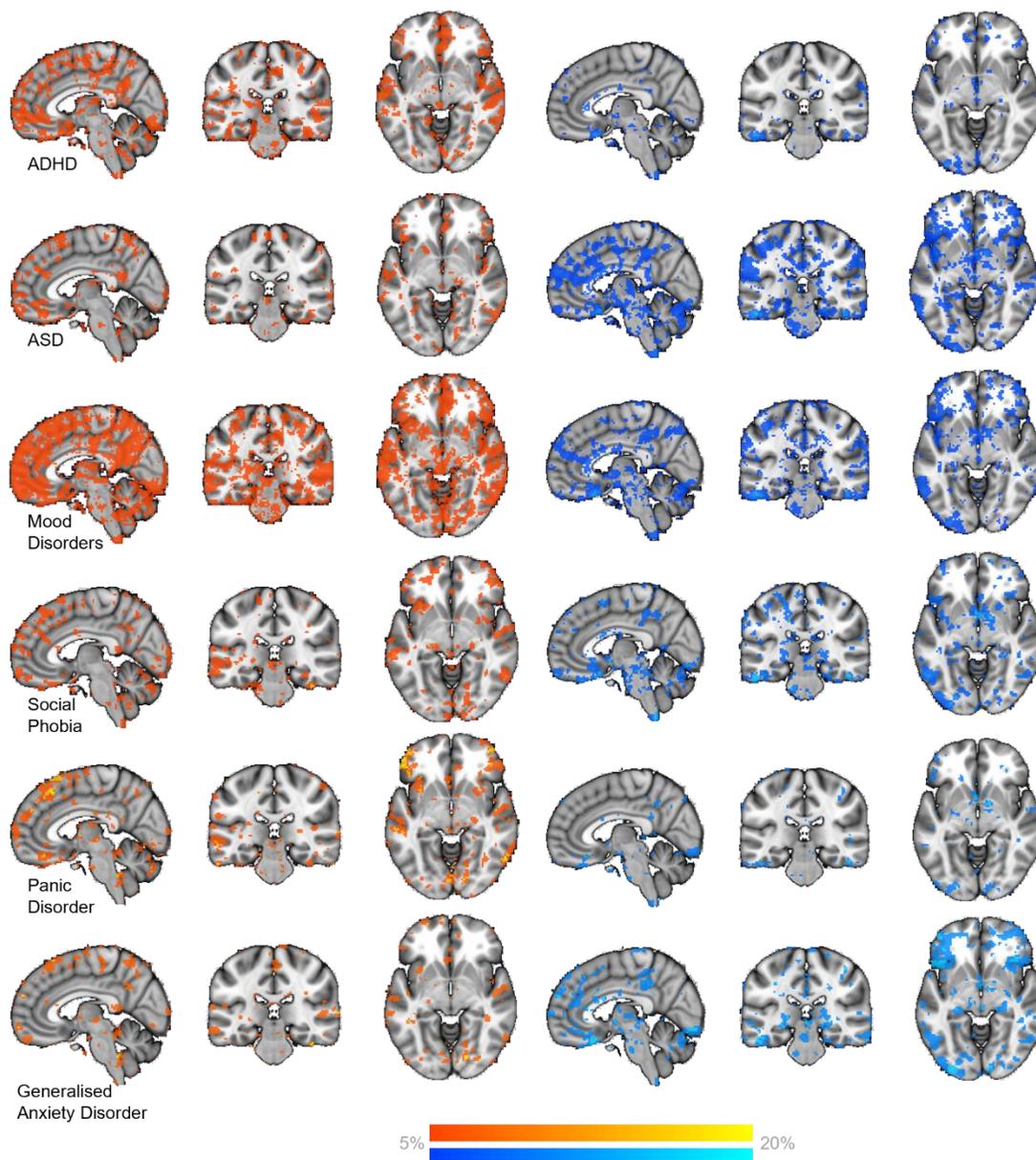
636 **Supplementary Figure 4: Evaluation of the faces>shapes (left) and faces>baseline normative models when applied to MIND-**
637 **Set cohort.** Histograms show the explained variance (a,d), skew (b,e), and kurtosis (c,f) of the clinical data, as tested on reference
638 normative models of EFMT related BOLD activation, and their respective illustration on the brain (x,y,z= 4,-6,-15).
639

640

641

642

643 **SUPPLEMENTARY RESULTS – Location of deviations for diagnoses:**



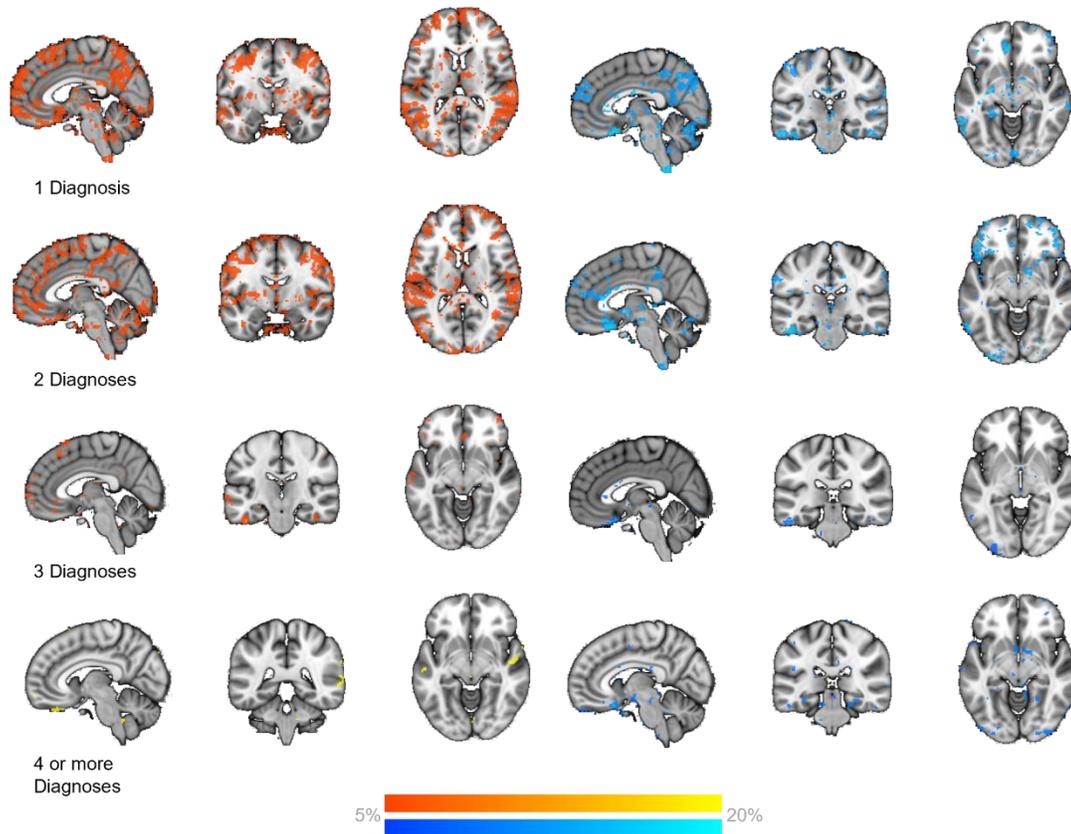
644
645
646
647
648
649
650
651

Supplementary Figure 5: Heterogeneous location of deviations in predicted BOLD signal for different types of neurodivergence, and mental health diagnoses. Maps illustrate the percentage of participants with a neurodivergence or mental health condition who had positive (left; hot colours) or negative deviations (right; cool colours) $> \pm 2.6$ within each voxel [minimum = %5 of sample, or 1 participant where 5% was a participant count less than 1, maximum = 20% of disorder sample size]. x,y,z, = 5, -28, -6.

652

653 **SUPPLEMENTARY RESULTS – Location of deviations for increasing levels of co-occurring**

654 **diagnoses:**



655

656

657 **Supplementary Figure 6: Heterogeneous location of deviations in predicted BOLD signal for increasing levels of co-**
658 **occurring diagnoses.** Maps illustrate the percentage of participants with a neurodivergence or mental health condition who had
659 positive (left; hot colours) or negative deviations (right; cool colours) $> \pm 2.6$ within each voxel [minimum = %5 of sample, or 1 participant
660 where 5% was a participant count less than 1, maximum = 20% of sample size].
661

662

663

664 **References:**

- 665 1 Marquand, A. F. *et al.* Conceptualizing mental disorders as deviations from normative
666 functioning. *Molecular Psychiatry* **24**, 1415-1424, doi:10.1038/s41380-019-0441-1 (2019).
- 667 2 Marquand, A. F., Rezek, I., Buitelaar, J. & Beckmann, C. F. Understanding Heterogeneity in
668 Clinical Cohorts Using Normative Models: Beyond Case-Control Studies. *Biological psychiatry* **80**,
669 552-561, doi:10.1016/j.biopsych.2015.12.023 (2016).
- 670 3 Rutherford, S. *et al.* Charting brain growth and aging at high spatial precision. *eLife* **11**, e72904,
671 doi:10.7554/eLife.72904 (2022).
- 672 4 Zabihi, M. *et al.* Fractionating autism based on neuroanatomical normative modeling.
673 *Translational Psychiatry* **10**, 384, doi:10.1038/s41398-020-01057-0 (2020).
- 674 5 Zabihi, M. *et al.* Dissecting the Heterogeneous Cortical Anatomy of Autism Spectrum Disorder
675 Using Normative Models. *Biological psychiatry. Cognitive neuroscience and neuroimaging* **4**,
676 567-578, doi:10.1016/j.bpsc.2018.11.013 (2019).
- 677 6 Bethlehem, R. A. I. *et al.* A normative modelling approach reveals age-atypical cortical thickness
678 in a subgroup of males with autism spectrum disorder. *Communications Biology* **3**, 486,
679 doi:10.1038/s42003-020-01212-9 (2020).
- 680 7 Wolfers, T. *et al.* Mapping the Heterogeneous Phenotype of Schizophrenia and Bipolar Disorder
681 Using Normative Models. *JAMA Psychiatry* **75**, 1146-1155,
682 doi:10.1001/jamapsychiatry.2018.2467 (2018).
- 683 8 Cropley, V. L. *et al.* Brain-Predicted Age Associates With Psychopathology Dimensions in Youths.
684 *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging* **6**, 410-419,
685 doi:<https://doi.org/10.1016/j.bpsc.2020.07.014> (2021).
- 686 9 Hariri, A. R., Tessitore, A., Mattay, V. S., Fera, F. & Weinberger, D. R. The amygdala response to
687 emotional stimuli: a comparison of faces and scenes. *Neuroimage* **17**, 317-323,
688 doi:10.1006/nimg.2002.1179 (2002).
- 689 10 Hariri, A. R. *et al.* Serotonin Transporter Genetic Variation and the Response of the Human
690 Amygdala. **297**, 400-403, doi:doi:10.1126/science.1071829 (2002).
- 691 11 Miller, K. L. *et al.* Multimodal population brain imaging in the UK Biobank prospective
692 epidemiological study. *Nat Neurosci* **19**, 1523-1536, doi:10.1038/nn.4393 (2016).
- 693 12 Van Essen, D. C. *et al.* The WU-Minn Human Connectome Project: an overview. *Neuroimage* **80**,
694 62-79, doi:10.1016/j.neuroimage.2013.05.041 (2013).
- 695 13 Van Essen, D. C. *et al.* The Human Connectome Project: a data acquisition perspective.
696 *Neuroimage* **62**, 2222-2231, doi:10.1016/j.neuroimage.2012.02.018 (2012).
- 697 14 Harms, M. P. *et al.* Extending the Human Connectome Project across ages: Imaging protocols for
698 the Lifespan Development and Aging projects. *Neuroimage* **183**, 972-984,
699 doi:10.1016/j.neuroimage.2018.09.060 (2018).
- 700 15 Snoek, L. *et al.* The Amsterdam Open MRI Collection, a set of multimodal MRI datasets for
701 individual difference analyses. *Scientific Data* **8**, 85, doi:10.1038/s41597-021-00870-6 (2021).
- 702 16 van Eijndhoven, P. *et al.* Measuring Integrated Novel Dimensions in Neurodevelopmental and
703 Stress-Related Mental Disorders (MIND-SET): Protocol for a Cross-sectional Comorbidity Study
704 From a Research Domain Criteria Perspective. *JMIRx Med* **3**, e31269, doi:10.2196/31269 (2022).
- 705 17 Mulders, P. C. R. *et al.* Striatal connectopic maps link to functional domains across psychiatric
706 disorders. *Translational Psychiatry* **12**, 513, doi:10.1038/s41398-022-02273-6 (2022).
- 707 18 Witten, D. M., Tibshirani, R. & Hastie, T. A penalized matrix decomposition, with applications to
708 sparse principal components and canonical correlation analysis. *Biostatistics* **10**, 515-534,
709 doi:10.1093/biostatistics/kxp008 (2009).

- 710 19 Witten, D. M. & Tibshirani, R. J. Extensions of Sparse Canonical Correlation Analysis with
711 Applications to Genomic Data. **8**, doi:doi:10.2202/1544-6115.1470 (2009).
- 712 20 Bayer, J. M. M. *et al.* Accommodating site variation in neuroimaging data using normative and
713 hierarchical Bayesian models. *NeuroImage* **264**, 119699,
714 doi:<https://doi.org/10.1016/j.neuroimage.2022.119699> (2022).
- 715 21 Nygaard, V., Rødland, E. A. & Hovig, E. Methods that remove batch effects while retaining group
716 differences may lead to exaggerated confidence in downstream analyses. *Biostatistics* **17**, 29-39,
717 doi:10.1093/biostatistics/kxv027 (2016).
- 718 22 Kebets, V. *et al.* Fronto-limbic neural variability as a transdiagnostic correlate of emotion
719 dysregulation. *Translational Psychiatry* **11**, 545, doi:10.1038/s41398-021-01666-3 (2021).
- 720 23 Rutherford, S. *et al.* Evidence for embracing normative modeling. *eLife* **12**, e85082,
721 doi:10.7554/eLife.85082 (2023).
- 722 24 Westlin, C. *et al.* Improving the study of brain-behavior relationships by revisiting basic
723 assumptions. *Trends in cognitive sciences*, doi:10.1016/j.tics.2022.12.015 (2023).
- 724 25 Everaerd, D., Klumpers, F., Oude Voshaar, R., Fernández, G. & Tendolkar, I. Acute Stress
725 Enhances Emotional Face Processing in the Aging Brain. *Biological Psychiatry: Cognitive
726 Neuroscience and Neuroimaging* **2**, 591-598, doi:<https://doi.org/10.1016/j.bpsc.2017.05.001>
727 (2017).
- 728 26 Mizzi, S., Pedersen, M., Lorenzetti, V., Heinrichs, M. & Labuschagne, I. Resting-state
729 neuroimaging in social anxiety disorder: a systematic review. *Molecular Psychiatry*,
730 doi:10.1038/s41380-021-01154-6 (2021).
- 731 27 Elliott, M. L. *et al.* What Is the Test-Retest Reliability of Common Task-Functional MRI Measures?
732 New Empirical Evidence and a Meta-Analysis. *Psychological Science* **31**, 792-806,
733 doi:10.1177/0956797620916786 (2020).
- 734 28 Somerville, L. H. *et al.* The Lifespan Human Connectome Project in Development: A large-scale
735 study of brain connectivity development in 5-21 year olds. *NeuroImage* **183**, 456-468,
736 doi:10.1016/j.neuroimage.2018.08.050 (2018).
- 737 29 Alfaro-Almagro, F. *et al.* Image processing and Quality Control for the first 10,000 brain imaging
738 datasets from UK Biobank. *NeuroImage* **166**, 400-424 (2018).
- 739 30 van Eijndhoven, P. F. P. *et al.* Measuring Integrated Novel Dimensions in Neurodevelopmental
740 and Stress-related Mental Disorders (MIND-Set): a cross-sectional comorbidity study from an
741 RDoC perspective. *medRxiv*, 2021.2006.2005.21256695, doi:10.1101/2021.06.05.21256695
742 (2021).
- 743 31 Oldehinkel, M. *et al.* Attention-Deficit/Hyperactivity Disorder symptoms coincide with altered
744 striatal connectivity. *Biological psychiatry. Cognitive neuroscience and neuroimaging* **1**, 353-363,
745 doi:10.1016/j.bpsc.2016.03.008 (2016).
- 746 32 Glasser, M. F. *et al.* The minimal preprocessing pipelines for the Human Connectome Project.
747 *NeuroImage* **80**, 105-124, doi:10.1016/j.neuroimage.2013.04.127 (2013).
- 748 33 Pruim, R. H. R. *et al.* ICA-AROMA: A robust ICA-based strategy for removing motion artifacts
749 from fMRI data. *NeuroImage* **112**, 267-277, doi:10.1016/j.neuroimage.2015.02.064 (2015).
- 750 34 Botvinik-Nezer, R. *et al.* Variability in the analysis of a single neuroimaging dataset by many
751 teams. *Nature* **582**, 84-88, doi:10.1038/s41586-020-2314-9 (2020).
- 752 35 Andersson, J. L., Jenkinson, M., Smith, S. & Oxford, F. A. G. o. t. U. o. Non-linear registration, aka
753 Spatial normalisation FMRIB technical report TR07JA2. **2**, e21 (2007).
- 754 36 Kraha, A., Turner, H., Nimon, K., Zientek, L. R. & Henson, R. K. Tools to support interpreting
755 multiple regression in the face of multicollinearity. *Frontiers in psychology* **3**, 44,
756 doi:10.3389/fpsyg.2012.00044 (2012).

